

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平10-95796

(43) 公開日 平成10年(1998) 4月14日

(51) Int.Cl.<sup>6</sup>

識別記号

F I

C 0 7 K 14/00

C 0 7 K 14/00

G 0 1 N 33/68

G 0 1 N 33/68

// G 0 6 F 17/00

G 0 6 F 15/20

D

審査請求 有 請求項の数 2 O L (全 8 頁)

(21) 出願番号 特願平9-283233  
(62) 分割の表示 特願平5-246805の分割  
(22) 出願日 平成5年(1993)10月1日

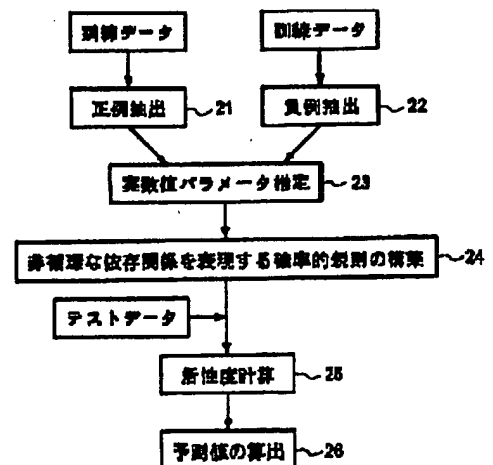
(71) 出願人 000004237  
日本電気株式会社  
東京都港区芝五丁目7番1号  
(72) 発明者 馬見塚 拓  
東京都港区芝五丁目7番1号 日本電気株  
式会社内  
(74) 代理人 弁理士 京本 直樹 (外2名)

(54) 【発明の名称】 タンパク質立体構造予測方法

(57) 【要約】

【課題】 構造未知のタンパク質のアミノ酸配列情報から、その立体構造の各局所構造を高精度で予測する。

【解決手段】 ステップ21で構造既知及び未知のタンパク質アミノ酸配列を入力とし、それらのアライメント(整合)から局所構造領域の正例を出力し、ステップ22で立体構造既知のタンパク質アミノ酸配列を入力とし、それらのアライメントから局所構造領域の負例を出力し、ステップ23で前記正例と負例を入力とし、これら訓練データのアミノ酸種類あるいは実数値属性から確率的規則の実数値パラメータの推定値を計算し、ステップ24では情報量規準に基づき、各残基位置に相当する変数間の依存関係を表現する確率的規則の構造を決定し、ステップ25でテストデータ配列を入力とし、前記確率的規則を使用し、テストデータ配列の各領域が局所構造であるかの確からしさを示す活性度を出力し、ステップ26で前記活性度を入力とし、その中から最適値を出力する。



## 【特許請求の範囲】

【請求項1】タンパク質構造の規則を学習するステップを有するタンパク質立体構造予測方法において、前記ステップが、正例と負例とからなる訓練データのアミノ酸配列を入力とし、各依存関係を条件付き確率とし、あらかじめ定められた情報量基準を用いて、該アミノ酸配列の各々の残基位置が依存している他の残基位置を、依存関係が循環しないという制限の下で決定することにより、該依存関係を表現する確率的規則を構築して出力することを特徴とするタンパク質立体構造予測方法。

【請求項2】タンパク質構造の規則を学習するステップを有するタンパク質立体構造予測方法において、前記ステップが、正例と負例とからなる訓練データのアミノ酸配列の実数値属性を入力とし、各依存関係を条件付き確率とし、あらかじめ定められた情報量基準を用いて、該アミノ酸配列の各々の残基位置が依存している他の残基位置を、依存関係が循環しないという制限の下で決定することにより、該依存関係を表現する確率的規則を構築して出力することを特徴とするタンパク質立体構造予測方法。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】本発明は、立体構造未知のタンパク質アミノ酸配列から、タンパク質の立体構造を予測する方法に関する。

## 【0002】

【従来の技術】タンパク質の局所的な立体構造として、 $\alpha$ ヘリックスや $\beta$ シートに代表される二次構造やジンクフィンガーやロイシンジッパーに代表されるモチーフなどがある。立体構造未知のタンパク質アミノ酸配列に対して、これらタンパク質の局所構造の予測が可能になれば、タンパク質全体の立体構造予測が可能になると一般に考えられている。

【0003】例えば、タンパク質二次構造予測問題は、20年以上も前から解決が図られてきた問題であり、従来、タンパク質の一次構造の各残基（以下、予測対象となる残基を中心残基と呼ぶ）が、3（あるいは4）種類の二次構造のいずれに対応するかを予測する問題として扱われてきた。従来技術によるタンパク質の二次構造を予測する方法として、例えば、1974年発行の米国の雑誌「バイオケミストリー」(Biochemistry)の第23巻222-245頁掲載のチョウ(Chou)とファスマン(Fasman)による論文「プレディクション オブ プロテイン コンホメーション」(Prediction of protein conformation)（以下、CF法と略す）、1978年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」(Journal of Molecular Biology)の第120巻97-120頁掲載のガルニエ(Garnier)らによる論文

「アナリシス オブ ザ アクキュレシー アンド インプリケーションズ オブ シンプル メソッド フォープレディクティング ザ セコンダリー ストラクチャー オブ グロブラープロテインズ」(Analysis of the accuracy and implications of simple methods for predicting the secondary structure of globular proteins)（以下、GOR法と略す）、1987年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」(Journal of Molecular Biology)の第198巻425-443頁掲載のギブラト(Gibrat)らによる論文「ファザー デベロップメント オブ プロテイン セコンダリー ストラクチャー プレディクションユージング インホメーション セオリー:ニュー パラメータズ アンド コンシダレーション オブ レジデュー ペアズ」(Further developments of protein secondary structure prediction using information theory: New parameters and consideration of residue pairs)（以下、GGR法と略す）、1988年発行の米国の雑誌「ジャーナル オブ モレキュラ バイオロジー」(Journal of Molecular Biology)の第202巻865-884頁掲載のキャン(Qian)らによる論文「プレディクティング ザ セコンダリー ストラクチャー オブ グロブラー プロテインズ ユージング ニューラル ネットワーク モデルズ」(Predicting the secondary structure of globular protein using neural network models)（以下QSと略す）、及び1993年の米国の学会「ハワイ インターナショナル コンファレンス オブ システム サイエンス」(Hawaii International Conference of System Sciences)の予稿集第1巻659-668頁記載の馬見塚らによる論文「プロテイン  $\alpha$ ヘリックス リージョン プレディクション ベースド オン ストキャスティックルール ラーニング」(Protein  $\alpha$ -Helix Region Prediction Based-on Stochastic Rule Learning)（以下、MY法と略す）などがある。

【0004】CF法は、タンパク質構造のデータベースから各二次構造におけるアミノ酸の統計的な出現頻度を求め、この頻度表を使用し、経験的な規則に基づく予測を行っている。また、GOR法は、中心残基の二次構造に対して、その残基から数残基離れた残基により独立に

もたらされる情報量の和を計算し、その相対値から予測を行い、GGR法は、中心残基の二次構造に対して、その残基及びその残基から数残基離れた残基によりもたらされる情報量の和から予測を行っている。さらに、QS法は、3層のフィードフォワード型のネットワークを使用し、中心残基の前後8残基を含む配列を入力とし、二次構造に対する中心残基及び周辺残基からの寄与をニューラルネットワークを用いて抽出することにより予測を行っている。MY法は、訓練配列の各残基位置において各アミノ酸が $\alpha$ ヘリックスであるかの確からしさを確率分布として計算し、それからテスト配列の各領域に対して、 $\alpha$ ヘリックスの確からしさを計算する。

#### 【0005】

【発明が解決しようとする課題】タンパク質アミノ酸配列においては、各アミノ酸残基同士は依存関係を保持し、その局所的な立体構造や機能的な部位を形成していると考えられている。従って、タンパク質の局所的な立体構造を表現及び予測するためには、それら局所的な立体構造内の各残基間の依存関係の表現が重要であると考えられる。しかし、従来、それら残基位置間に依存する依存関係をネットワークの形で自動的に抽出する方法や、さらに、その依存関係を規則として未知データに対する予測を行う方法は、全く検討されておらず、そういった手法が確率されていなかった。

#### 【0006】

【課題を解決するための手段】第1の発明は、タンパク質構造の規則を学習するステップを有するタンパク質立体構造予測方法において、前記ステップが、正例と負例とからなる訓練データのアミノ酸配列を入力とし、各依存関係を条件付き確率とし、あらかじめ定められた情報量基準を用いて、該アミノ酸配列の各々の残基位置が依存している他の残基位置を、依存関係が循環しないという制限の下で決定することにより、該依存関係を表現する確率的規則を構築して出力することを特徴とする。

【0007】第2の発明は、タンパク質構造の規則を学習するステップを有するタンパク質立体構造予測方法において、前記ステップが、正例と負例とからなる訓練データのアミノ酸配列の実数値属性を入力とし、各依存関係を条件付き確率とし、あらかじめ定められた情報量基準を用いて、該アミノ酸配列の各々の残基位置が依存している他の残基位置を、依存関係が循環しないという制限の下で決定することにより、該依存関係を表現する確率的規則を構築して出力することを特徴とする。

#### 【0008】

【発明の属する技術分野】次に、本発明について図面を参照して詳細に説明する。

【0009】図1は、本発明に関連するタンパク質立体構造予測方法の実施例を説明するフローチャートである。本実施例では、対象とするタンパク質の局所的な立体構造として $\alpha$ ヘリックスを扱うものとする。

【0010】ステップ11では、 $\alpha$ ヘリックスの領域がわかっているタンパク質のアミノ酸配列に対して、同じファミリーのタンパク質、例えば、種が異なる同じタンパク質のアライメント（整合）をとり、 $\alpha$ ヘリックスに対応する部分配列を、 $\alpha$ ヘリックスの正例として抽出する。

【0011】例えば、ヘモグロビンというタンパク質の $\beta$ 鎖の場合には、ヒトのヘモグロビンの $\alpha$ ヘリックスの位置は、X線結晶回折の結果から明らかになっており、8個の $\alpha$ ヘリックスの領域を有することが知られている。従って、ヒトのヘモグロビン $\beta$ 鎖に対して、他の種、例えば、チンパンジー、ウマなどの他の種のヘモグロビン $\beta$ 鎖のアライメントを行い、8個の $\alpha$ ヘリックスに対応する領域を $\alpha$ ヘリックスの正例として抽出する。

【0012】ステップ12では、 $\alpha$ ヘリックス位置の知られているタンパク質の $\alpha$ ヘリックスに対応する部分配列に対して、 $\alpha$ ヘリックス位置の知られているアミノ酸配列データベースの各配列のアライメントをとり、 $\alpha$ ヘリックスに対応しない部分配列を、ステップ10で抽出された $\alpha$ ヘリックスの正例に対する負例として抽出する。

【0013】ヘモグロビン $\beta$ 鎖の例では、8個の $\alpha$ ヘリックスに対応する部分配列に対して、例えば、PDB (Protein Data Bank) などのタンパク質構造データベース内のいくつかの配列に対してアライメントを行い、アライメントの結果得られた各部分配列において、その配列の構造が $\alpha$ ヘリックスではない場合に、それらを負例として抽出する。例えば、負例抽出の際のアライメントでは、一定の割合以上の相同性を保持する部分配列を負例とすることが考えられる。具体的には、アライメントによる相同性が30%以上の部分配列を負例とする方法などがある。

【0014】抽出するデータ数については、例えば、 $\alpha$ ヘリックスの正例となる各領域における正例と負例との割合を各領域についてそれぞれ等しくすることが考えられる。さらに、具体的には、その割合として正例、負例を同数とすることが考えられる。

【0015】ステップ13は、ステップ11及びステップ12で抽出された正例と負例から、確率的規則の実数値パラメータを計算するステップである。

【0016】確率的規則とは、ここでは任意の与えられた配列の領域に対して、 $\alpha$ ヘリックスに対応する確率を与える確率分布のことである。各 $X_i$  ( $i = 1, \dots, n$ ) をそれぞれ属性値の空間として、 $X$  をそれらの直積、すなわち、 $X = X_1 \times X_2 \times \dots \times X_n$  と書く。

【0017】例えば、 $X$ は20種類のアミノ酸からなる一つの集合を表す場合や、また $X = X_1 \times X_2$  で、 $X_1$  が疎水性を表す数値の範囲かつ $X_2$  が分子量を表す数値の範囲を表す場合などがある。

【0018】 $\alpha$ ヘリックスの正例中の長さLのウィンドウWに対し、テスト配列中の任意の長さLの領域SがW部分に対応する確からしさを以下のように求める。まず、 $X_t$  (以下、変数と呼ぶ)を配列Sの左から数えてt番目の残基位置であり、 $\pi_t$ を領域Wにおいて $X_t$ が依存する残基位置の集合(以下、 $\pi_t$ を $X_t$ の親変数の集合、 $X_t$ を $\pi_t$ の子変数と呼ぶ)とする。ここで、

【0019】

\*【数1】

$$P(X_t|\pi_t)$$

【0020】を、領域Wにおいて、 $X_t$ が $\pi_t$ に依存している条件付き確率とし、領域SがW部分に対応する確からしさ $P_W(S)$ は、次のように書けるものと仮定する。

【0021】

\*【数2】

$$P_W(S) = \prod_{t=1}^L P_W(X_t|\pi_t) \quad (1)$$

【0022】1式の右辺は、変数をノードとし、変数間の親から子にアークを伸ばすことにより、ネットワーク構造に対応する。例えば、領域Sが3残基からなり、領域Sの各残基の結合確率が具体的に次式のように書ける※

※ものとすれば、次式は図3のネットワークに対応する。

【0023】

【数3】

$$P_W(S) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \quad (2)$$

【0024】さらに、各

【0025】

【数4】

$$P(X_t|\pi_t)$$

【0026】は、与えられた正例と負例とからなる事例データから、例えば、次のようにして決定される。

【0027】まず、t番目の残基位置において、属性の実数値のとり得る範囲を重なり合わない部分領域(以下、これをセルと呼ぶ)に有限分割し、 $m_t$ を全セル数、 $C_i$ をi番目のセルとする。

【0028】t番目の位置の残基が $m_t$ 個のセルの中の $C_i$ に含まれる場合に、 $X_t$ の生起確率 $P(X_t = i) = p_i(t)$ とする。ここで、

【0029】

【数5】

$$p_i(t) \in [0, 1] \quad (i = 1, \dots, m_t)$$

20★【0030】であり、これを確率パラメータと呼ぶ。図4は、有限分割の構造を示す例であるが、値が0から1の範囲をとる一つの属性により確率パラメータを推定する場合を示す。

【0031】確率パラメータは、各セルに含まれる正例及び負例のデータ数を用いて推定する。

【0032】

【数6】

$$N_i^+(t)$$

30【0033】をt番目の位置でのi番目のセルに含まれる正例数、 $N_i^-(t)$ をt番目の位置でのi番目のセルに含まれる正例数と負例数の和とし、t番目の位置でのi番目のセルにおける推定値を $p_i(t)$ とする。例えば、次式のラプラス推定量によって、各セルに対する確率パラメータを計算する。

【0034】

★【数7】

$$p_i(t) = \frac{N_i^+(t) + 1}{N_i(t) + 2} \quad (3)$$

$$(i = 1, \dots, m_t)$$

【0035】ただし、推定量はラプラス推定量のみならず、多くの推定量が使用できる。次に、同様に、 $X_t$ と $\pi_t$ の結合確率 $P(X_t, \pi_t)$ も推定量を用いて計算できる。例えば、 $\pi_t$ の要素が変数 $X_s$ のみである場合、 $X_s$ を $X_t$ と同様に重なり合わない $m_s$  ☆

☆個の部分領域に有限分割し、s番目の残基が $m_s$ 個のセル内の $C_j$ に含まれ、

【0036】

【数8】

$$P(X_i = i, X_s = j) = p_{ij}(t, s) \quad (p_{ij}(t, s) \in [0, 1] \quad (i = 1, \dots, m_t, j = 1, \dots, m_s))$$

【0037】とし、確率パラメータ  $p_{ij}(t, s)$  を推定する。

【0038】  $t$  番目、  $s$  番目の各位置において、各セルに含まれる正例及び負例のデータ数から、

【0039】

【数9】

$$N_{ij}^+(t, s)$$

\*

$$p_{ij}(t, s) = \frac{N_{ij}^+(t, s) + 1}{N_{ij}(t, s) + 2} \quad (4)$$

$$(i = 1, \dots, m_t, j = 1, \dots, m_s)$$

【0042】最後に、これら推定された確率パラメータを用いて、  $\pi_t$  が存在する下での  $X_t$  の条件付き確率

【0043】

【数11】

$$P(X_t | \pi_t)$$

\* 【0040】を各位置においてそれぞれ  $i, j$  番目のセルに含まれる正例数、  $N_{ij}(t, s)$  を各位置においてそれぞれ  $i, j$  番目のセルに含まれる正例数と負例数の和とする。これから、例えば、次式のラプラス推定量により、確率パラメータを推定する。

【0041】

【数10】

20

【0045】

【数12】

【0044】を確率パラメータとして計算する。上述の※

$$P(X_i = i | X_j = j) = p_{ij}(t, s) \quad (p_{ij}(t, s) \in [0, 1] (i = 1, \dots, m_t, j = 1, \dots, m_s))$$

【0046】とし、確率パラメータ

【0047】

【数13】

★ 【0048】は次式のように計算する。

【0049】

【数14】

$$p_{ij}(t | s)$$

★

$$p_{ij}(t | s) = \frac{p_{ij}(t, s)}{p_j(s)} \quad (5)$$

$$(i = 1, \dots, m_t, j = 1, \dots, m_s)$$

【0050】ステップ14では、確率的規則の構造を決定する。すなわち、各変数  $X_t$  に対し、その親変数の集合  $\pi_t$  を情報量基準を使用して決定するステップである。

【0051】以下、情報量基準として記述長最小 (Minimum Description Length (MDL)) 原理 (以下、MDL原理) と適用した場合のネットワーク構成方法の一例について述べる。なお、MDL原理については、1978年発行の米国の雑誌

「オートマティカ」 (Automatica) の第14巻465-471頁記載のリサネン (Rissanen) による論文「モデリングバイ ショーテスト データ ディスクリプション」 (Modeling by shortest data description) に詳しく記載されている。

【0052】MDL原理によれば、与えられた事例デー

タから計算されるデータ記述長と規則の記述長との和が最小となる規則を最適な規則とする。従って、ステップ11、12、13において求められた正例、負例及び実数値パラメータから、ここでの確率的規則のデータ記述長及び規則の記述長を計算する。

【0053】ここで説明する例では、各残基位置ごとに、その位置と依存関係にある位置を決定していくことを考える。すなわち、各変数毎に独立にその親の変数を決定していく。

【0054】まず、  $t$  番目の残基位置に着目する。変数  $X_t$  とその親変数の集合  $\pi_t$  との依存関係は、1式の確率的規則から条件付き確率

【0055】

【数15】

$$P(X_t | \pi_t)$$

【0056】で表現される。ここで、親変数の数を $k$ 、親変数の残基位置を順番に $t_1$  から  $t_k$  とし、また、 $t$  及び  $t_1$  から  $t_k$  の残基位置での全セル数をそれぞれ  $m_t, m_{t_1}, \dots, m_{t_k}$  とする。さらに、変数  $X_t$  の残基が  $i$ 、変数  $X_{t_1}$  から  $X_{t_k}$  の残基が、それぞれ  $j_1, \dots, j_k$  番目のセルに含まれるような属性を有している正例数を

【0057】

【数16】

$$N_{i,j_1,\dots,j_k}^+(t, t_1, \dots, t_k)$$

【0058】、変数  $X_t$  の残基が  $i$ 、変数  $X_{t_1}$  から  $X_{t_k}$  の残基が、それぞれ  $j_1, \dots, j_k$  番目のセルに含まれるような属性を有している正例数と負例数との和を  $N_{i,j_1,\dots,j_k}(t, t_1, \dots, t_k)$

$$- \sum_{j_1=0}^{m_{t_1}} \dots \sum_{j_k=0}^{m_{t_k}} \sum_{i=0}^{m_t} \log(p_{i,j_1,\dots,j_k}(t|t_1, \dots, t_k) N_{i,j_1,\dots,j_k}^+(t, t_1, \dots, t_k)) \\ \times (1 - p_{i,j_1,\dots,j_k}(t|t_1, \dots, t_k))^{N_{i,j_1,\dots,j_k}^-(t, t_1, \dots, t_k) - N_{i,j_1,\dots,j_k}^+(t, t_1, \dots, t_k)} \quad (6)$$

【0063】さらに、ここでの確率的規則の規則の記述長は、次式で与えられる。

※【0064】

※【数19】

$$- \sum_{j_1=0}^{m_{t_1}} \dots \sum_{j_k=0}^{m_{t_k}} \sum_{i=0}^{m_t} \frac{\log N_{i,j_1,\dots,j_k}(t, t_1, \dots, t_k)}{2} \quad (7)$$

【0065】従って、 $t$  番目の残基位置に相当する変数  $X_t$  に対し、(6) 式と (7) 式との和を最小にするような親変数の集合を選択することにより、確率的規則の構造が決定される。

【0066】ステップ15では、ステップ14において構成された確率的規則を使用し、与えられたテストデータ配列の各領域に対して、その活性度を計算する。

【0067】ここでは、活性度として尤度を使用する。

【0068】まず、訓練配列の正例中の任意の長さ  $L$  の領域  $W$  を取り出す。この  $W$  の各残基位置に対応する変数に対して、その親変数がステップ14において決定されている。さらに、各変数とその親変数との依存関係を表す条件付き確率の実数値パラメータは、ステップ13に★40

$$P_W(S) = p_2(1)p_{1,2}(2|1)p_{3,2,1}(3|1,2) \quad (8)$$

【0072】この動作を訓練配列の正例中の取り得る全ての長さ  $L$  の領域で構成された確率的規則を使用し、テスト配列中の長さ  $L$  の全ての部分配列に対して行う。

【0073】ステップ16では、テスト配列中の任意の長さ  $L$  の領域  $S$  に対して、訓練配列中の正例の取り得る全ての長さ  $L$  の領域により算出された複数の尤度の中で、最大の尤度を選出し、領域  $S$  の  $\alpha$  ヘルックスの尤度とする。

\*、 $t_k$  )、変数  $X_t$  の残基が  $i$ 、変数  $X_{t_1}$  から  $X_{t_k}$  の残基が、それぞれ  $j_1, \dots, j_k$  番目のセルに入るような属性を有している条件付き確率の確率パラメータを

【0059】

【数17】

$$p_{i,j_1,\dots,j_k}(t|t_1, \dots, t_k)$$

【0060】とする。

10 【0061】すると、ここでの確率的規則のデータ記述長は、ステップ13により計算された確率パラメータから、規則の対数尤度の負をとることにより、次式で与えられる。

【0062】

【数18】

★おいて算出されている。

30 【0069】次に、この領域  $W$  をテスト配列の任意の長さ  $L$  の部分配列  $S$  にあてはめ、その  $\alpha$  ヘルックスの尤度を計算する。

【0070】例えば、 $L=3$  の領域  $W$  において、(2) 式のような確率的規則の構造が決定され、テスト配列の  $L=3$  の領域  $S$  では、その領域内の各残基が順に、2、1、3 番目のセルに入る実数値属性を有しているとすると、この領域  $S$  が  $W$  である尤度は次式で計算できる。

【0071】

【数20】

【0074】ステップ15及びステップ16の動作は次のようにまとめることができる。すなわち、訓練配列中の正例中の長さ  $L$  の領域の全ての集合を  $A$  とし、テスト配列の長さ  $L$  の部分配列  $S$  に対する  $\alpha$  ヘルックスの尤度  $P(S)$  を次式により計算する。

【0075】

【数21】

$$P(S) = \max_W^A P_W(S) \quad (9)$$

【0076】この動作をテスト配列の各領域に対して繰り返すことにより、テスト配列の各領域において、 $\alpha$ ヘリックスの尤度を計算する。

【0077】ここで、さらに、 $\alpha$ ヘリックス領域が複数個あれば、各 $\alpha$ ヘリックス領域について、同様な尤度計算を行ない、 $\alpha$ ヘリックス領域全体を通じて最大の尤度を最適値として選ぶ。

【0078】さらに、テスト配列内の尤度が与えられた各領域において、最大の尤度を領域内の各残基の最適値とする、あるいは、領域内の各残基に対しては、その残基を含む領域の得られた尤度の平均を各残基の最適値とする、などの方法を使用し、テストアミノ酸配列全体に対する尤度の変化を出力する。

【0079】以上の図1における学習及び予測方法は、 $\alpha$ ヘリックス以外の二次構造及びモチーフ等の局所領域の特徴抽出、さらに予測に適用できる。図2は、本発明のタンパク質立体構造予測方法の実施例を説明するフローチャートである。本実施例では、対象とする二次構造として $\alpha$ ヘリックスを扱うものとする。

【0080】ステップ21は、図1のステップ11と同じ処理を行ない $\alpha$ ヘリックス領域予測のために必要な正例を抽出する。

【0081】ステップ22は、図1のステップ12と同じ処理を行ない $\alpha$ ヘリックス領域予測のために必要な負例を抽出する。

【0082】ステップ23は、図1のステップ13と同じ処理を行ないステップ21及びステップ22で抽出された正例及び負例から、確率的規則の実数値パラメータを推定する。

【0083】ステップ24は、確率的規則の構築に制限が加えられ、局所構造領域の全変数の結合確率分布として無矛盾な確率的規則を構築するステップであり、本発明の第1の発明と第2の発明に含まれる。ここでの制限とは、確率的規則を図3のようなネットワーク構造で示した場合に、確率分布に矛盾が生じないように、アークの方向を非循環とする制限である。例えば、図3は非循環ネットワークの例であるが、この図において、 $X_1$  から $X_3$  に伸びているアークを逆に $X_3$  から $X_1$  へと伸ばせば、このネットワークは循環ネットワークとなり、そのようなネットワークの生成は許さない。

【0084】制限を加える方法として、例えば、各変数に順番付けを行ない、順番の小さい変数のみを親変数として持てるとする方法、あるいは、アークに循環が生じるとする依存関係が構成される場合のみ、その依存関

係を成立しないようにする方法などが考えられる。

【0085】ステップ25は、図1のステップ15と同じ処理を行ない、ステップ24を使用して構造が最適化された確率的規則を使用し、テストアミノ酸配列データの各領域に対して、その活性度を計算する。

【0086】ステップ26は、図1のステップ16と同じ処理を行ないステップ25により求められた複数の活性度から、配列全体に対する活性度の変化を出力する。

【0087】以上の図2における学習及び予測方法は、 $\alpha$ ヘリックス以外の二次構造及びモチーフ等の局所領域の特徴抽出、さらに予測に適用できる。

【0088】

【発明の効果】立体構造既知のタンパク質のアミノ酸配列情報から、局所的な立体構造さえも未知のタンパク質の局所的な立体構造を従来技術に対して高い精度で予測可能である。例えば、従来手法の一つであるMY法では、局所領域内の残基位置間の依存関係を全く考慮していなかったが、残基位置間の依存性を反映した確率的規則の構成によって、より精度の高い局所領域の特徴抽出及び予測が可能になっている。また、情報量規準に基づく最適化により、確率的規則の構造を理論的に最適化することが可能になる。

【図面の簡単な説明】

【図1】本発明に関連するタンパク質立体構造予測方法の一実施例を示すフローチャート

【図2】本発明のタンパク質立体構造予測方法の一実施例を示すフローチャート

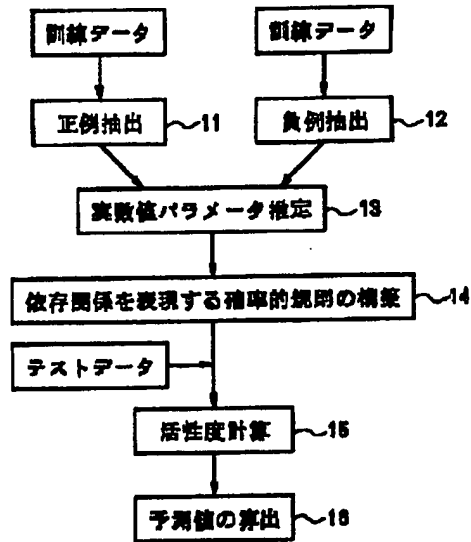
【図3】本発明で使用する確率的規則の変数間の依存関係を示す模式図。

【図4】本発明において各残基位置で行う有限分割の具体を示す模式図

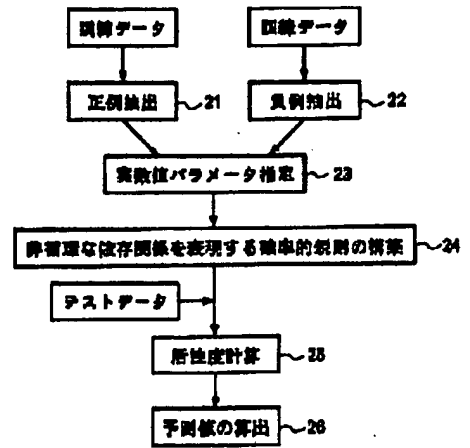
【符号の説明】

- 1 1 正例抽出
- 1 2 負例抽出
- 1 3 実数値パラメータ推定
- 1 4 確率的規則の構造の決定
- 1 5 テスト配列に対する活性度算出
- 1 6 テスト配列に対する予測値算出
- 2 1 正例抽出
- 2 2 負例抽出
- 2 3 実数値パラメータ推定
- 2 4 確率的規則の構造の決定
- 2 5 テスト配列に対する活性度計算
- 2 6 テスト配列に対する予測値算出

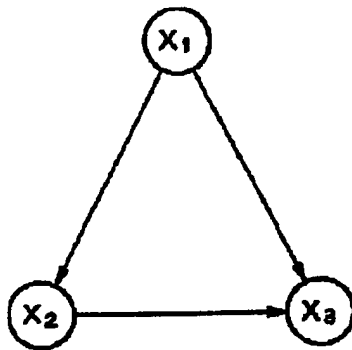
【図1】



【図2】



【図3】



【図4】

